

#### UNITED STATES PATENT APPLICATION

of

Sanjay Anand

Kyle Brandon

Nk Srinivas

and

Jameel Hyder

for

# METHOD AND COMPUTER PROGRAM PRODUCT FOR OFFLOADING PROCESSING TASKS FROM SOFTWARE TO HARDWARE

#### BACKGROUND OF THE INVENTION

#### 1. Cross-Reference to Related Applications

The present application is a continuation-in-part application of United States patent application serial number 09/657,510 filed 7 September 2000 and entitled "METHOD AND COMPUTER PROGRAM PRODUCT FOR OFFLOADING PROCESSING TASKS FROM SOFTWARE TO HARDWARE" which application is incorporated herein by reference in its entirety.

PH

#### 2. The Field of the Invention

The present invention relates generally to methods for increasing the efficiency, speed and/or throughput of a computer system. More specifically, the invention relates to methods for offloading computing tasks that are typically performed by a host processor in software, to a specific hardware component, thereby freeing up host processor resources and increasing the overall efficiency of the computer system.

#### 3. The State of the Art

A functional computer system generally consists of three fundamental components. The first component is the host computer and its associated peripheral hardware components. The host computer typically includes a central processing unit (CPU), which is interconnected via a bus with, for instance, system memory such as RAM or ROM. A system will also include a number of peripheral hardware devices, depending on the functionality needed, such as magnetic or optical disk storage devices, a keyboard or other input device, a display or other output device and communication equipment, such as a modem and/or a network interface card (NIC). Another fundamental computer component

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

1, 2

The street of th

ij

10 40

1214

1,2 

is the application software. Such software includes the familiar word processing applications, spread sheet applications, database applications, communications and network applications and so forth.

The final component of a modern, functional computer system is an operating system. The computer operating system performs many functions such as allowing a user to initiate execution of an application program. In addition, modern operating systems also provide an interface between application software and the host computer and its peripheral hardware. Thus, while it was once commonplace for an application program to directly access computer system hardware, modern operating systems provide standardized, consistent interfaces that allow user applications to interface with or access computer hardware peripherals in a standardized manner. To provide a consistent interface, operating system architectures are increasingly designed so that there may be several software layers between the actual hardware peripheral and the application program. For example, an application may make a call into the operating system. The operating system, in turn, may utilize the services provided by a hardware device driver layer. The device driver layer would then interface directly with the specific hardware peripheral. A primary advantage of such a layered approach is that layers may be added or replaced without impacting the other layers.

As will be appreciated, the complexity and sophistication of such operating systems, application software, and networking and communications continues to increase. This of course results in more functional and useful computer systems. However, this increased functionality is not without a cost. More feature rich operating systems and software applications often result in an increase in the processor overhead as a result of the additional duties that must be performed by a processor/CPU when executing such system

[]

60 EAST SOUTH TEMPLE SALT LAKE CITY, UTAH 84111 functions and/or applications. This phenomenon is especially apparent in connection with particular types of applications, such as network communication-type software applications. With the high bandwidth media that is increasingly prevalent, network speeds often match or exceed the CPU processor speed and memory bandwidth of the host computer. As such, to efficiently communicate over such networks, the CPU utilization and memory bandwidth used of the network-connected host computer must be minimized.

In addition, network applications further burden the host processor due to the layered architecture used by most, such as the seven-layer ISO model, or the layered model used by the Windows NT operating system. As is well known, such a model is used to describe the flow of data between the physical connection to the network and the end-user application. The most basic functions, such as putting data bits onto the network cable, are performed at the bottom layers, while functions attending to the details of applications are at the top layers. Essentially, the purpose of each layer is to provide services to the next higher layer, shielding the higher layer from the details of how services are actually implemented. The layers are abstracted in such a way that each layer believes it is communicating with the same layer on the other computer that is being communicated with via the network.

As will be appreciated, the various functions that are performed on a data packet as it proceeds between layers can be software intensive, and thus can demand a substantial amount of CPU processor and memory resources. For instance, in the Windows NT networking model, certain functions that are performed on the packet at various layers are extremely CPU intensive, such as packet checksum calculation and verification; encryption and decryption of data (e.g., SSL encryption and IP Security encryption); message digest calculation, TCP segmentation, receive side packet classification, packet filtering to guard

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

The state was seen that the state with the state with the state was seen to see the state with t

la de | d | 12 14

60 EAST SOUTH TEMPLE SALT LAKE CITY, UTAH 84111

against denial of service attacks, and User Datagram Protocol (UDP) send side packet fragmentation. As each of these functions are performed, the resulting demands on the CPU/memory can greatly effect the throughput and performance of the overall computer system.

Although software applications and operating system functions are placing greater demands on computer system resources, at the same time the capability, efficiency, and throughput of many computer hardware peripherals -- such as network interface cards (NICs) -- is also increasing. These computer system peripherals are often equipped with a dedicated processor and memory, and typically are capable of performing very sophisticated and complex computing tasks -- tasks that are otherwise performed by the computer system processor in software. For instance, many NICs are capable of independently performing tasks otherwise performed by the CPU in software at an checksum calculation/verification; data layer, such as appropriate network encryption/decryption; message digest calculation; TCP or UDP segmentation; receive side packet classification; packet filtering to guard against denial of service attacks; and others. As such, there is an advantage in offloading such CPU intensive task to a peripheral hardware device. This would reduce processor utilization and memory bandwidth usage in the host computer, and thereby increase the efficiency, speed and throughput of the overall system.

However, the processing capabilities of different peripheral devices vary widely. Thus, there needs to be an efficient method by which a computer system/operating system can identify the processing capabilities of such peripheral devices, and then assign and offload specific processing tasks to the device when needed. Also, it would be desirable if the tasks could be identified and assigned dynamically, depending on the then current needs of the processor. This would allow the computer system processor to take advantage of the capabilities of a hardware peripheral on an as-needed basis.

The state of the s

1.4

## 

. .

#### SUMMARY AND OBJECTS OF THE INVENTION

The foregoing problems in the prior state of the art have been successfully overcome by the present invention, which is directed to a system and method for offloading functions and tasks that were previously performed at a processor-software level, to an appropriate hardware peripheral connected to the computer system. The invention is particularly useful in connection with the offloading of tasks to network interface card (NIC) peripheral devices, which can often perform many of the tasks otherwise performed by the computer CPU in software.

In a preferred embodiment of the invention, a software implemented method and protocol is provided that allows, for instance, the operating system (OS) to "query" the device drivers (often referred to as "MAC" drivers) of any hardware peripherals (such as a NIC) that are connected to the computer system. The various device drivers each respond by identifying their respective hardware peripheral's processing capabilities, referred to herein as "task offload capabilities." In the preferred embodiment, once the task offload capabilities of each particular peripheral have been identified, the OS can then enable selected peripherals to perform certain tasks that could potentially be used by the OS. The OS can thereafter request that a peripheral perform the previously enabled task, or tasks, in a dynamic, as-needed basis, depending on the then current processing needs of the computer system.

While this general inventive concept would be applicable to other application or operating system environments, embodiments of the current invention are described herein as being implemented and utilized in connection with the layered networking model of Windows NT. Of course, the invention could be implemented in connection with essentially any similar type of architecture for managing and controlling network

WORKMAN, NYDEGGER & SEELEY A PROFESSIONAL CORPORATION

PL

Harry H.

communications. Specifically, the invention provides the ability to offload tasks or functions that are typically performed on a network packet at, for instance, the various network layers, and which typically require dedicated CPU and memory resources. These offloaded tasks can instead be optionally performed by the hardware peripheral that provides the actual physical communications channel to the network -- the NIC. For instance, rather than perform certain of the CPU intensive operations on the data packet as it passes through the respective network layers -- e.g. checksum calculation/verification, encryption/decryption, message digest calculation; TCP or UDP segmentation, receive side packet classification; packet filtering to guard against denial of service attacks; -- those tasks can instead be offloaded and performed at the NIC hardware.

In a preferred embodiment of the present invention, in the Windows NT layered networking architecture, a transport protocol driver, or transport, is implemented with an appropriate program method so as to be capable of querying each of the device driver(s) associated with the corresponding NIC(s) connected to the computer. Each queried device driver is similarly implemented so as to be capable of responding by identifying its specific processing, or "task offload" capabilities. In a preferred embodiment, once the task offload capabilities of each individual peripheral device have been identified, the transport sets which of those specific capabilities are to be enabled. This essentially informs the peripheral device what type of tasks it should expect to perform during subsequent transmissions and/or receptions of data packets. Thereafter, the transport is able to take advantage of the enabled capabilities of a peripheral device on an as-needed basis. Preferably, the enabled functions are invoked via appropriate data that is appended to the actual data packet destined for the network channel. In this way, tasks can be offloaded dynamically, and more than one task can be offloaded at a time.

WORKMAN, NYDEGGER & SEELEY

Thus, before a network packet is to be sent to a particular lower level device driver (e.g., residing at the MAC sublayer in a Windows NT environment), the transport will first determine what the capabilities of the corresponding NIC are. If capable of a particular function or functions, the transport enables the desired functions. If during subsequent packet transmissions the transport desires that a particular task be offloaded to hardware, it can dynamically append information to the packet that signifies that the desired function(s) should be performed on that packet at the NIC hardware. For instance, the transport will set a data flag in the data packet, thereby notifying the corresponding device driver that the NIC should calculate and append a checksum to that outgoing packet. hardware/software on the corresponding NIC will then handle this particular packet processing on its own, without any intervention or assistance from the system CPU. The system processor is thus freed up to perform other processing tasks, and the overall efficiency and throughput of the system is improved.

As noted, in a preferred embodiment, tasks are downloaded dynamically. That is, the capability of the NIC can be selectively used on a per-packet basis, depending on the then current needs of the computer system. Moreover, since tasks that were previously performed at various levels of the network stack are now performed at a single point -- the NIC itself -- the approach is more tightly integrated and efficient, further improving the throughput of the entire system. Preferably, embodiments of the current invention provide the transport with the ability to "batch" operations, i.e., offload multiple tasks to a single NIC. For instance, a single NIC can perform both checksumming, and encryption on a packet, thereby eliminating multiple CPU cycles that would have otherwise been needed if the same functions were implemented at the respective network layers in software.

24

1

2

3

4

5

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

4.

STATE STATE

12 1

ļ. 25 [a de

FLE

2

3

4

5

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

in dis

Accordingly, it is a primary object of this invention to provide a system and method for offloading computing tasks from a computer system processor to a hardware peripheral connected to the computer system. It is another object of this invention to provide a system and method for identifying the processing capabilities of individual peripherals. It is still another object of the present invention to provide a system and method for offloading tasks that can be effectively used in connection with a layered network architecture, whereby tasks that are typically performed at various layers in the network are instead offloaded to the appropriate network interface card (NIC). A related object of the invention is to provide a system and method in which computing tasks can be offloaded in a dynamic, as-needed basis, depending on the then current processing state of the computer processor. Yet another object of the present invention is to provide a system and method in which multiple tasks can be batched together, and then offloaded to a single peripheral device, such as a NIC.

Additional objects and advantages of the invention will be set forth in the description which follows, and in part will be obvious from the description, or may be learned by the practice of the invention. The objects and advantages of the invention may be realized and obtained by means of the instruments and combinations particularly pointed out in the appended claims. These and other objects and features of the present invention will become more fully apparent from the following description and appended claims, or may be learned by the practice of the invention as set forth hereinafter.

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

1 2

1, 2 ų.į

;; =; i, =; 

1 22

127

#### BRIEF DESCRIPTION OF THE DRAWINGS

In order that the manner in which the above-recited and other advantages and objects of the invention are obtained, a more particular description of the invention briefly described above will be rendered by reference to specific embodiments thereof which are illustrated in the appended drawings. Understanding that these drawings depict only typical embodiments of the invention and are not, therefore, to be considered to be limiting of its scope, the invention will be described and explained with additional specificity and detail through the use of the accompanying drawings in which:

Figure 1 is a diagram representing an example of a typical computer system and attached peripheral devices that can be used in connection with the present invention;

Figure 2 is a diagram illustrating some of the functional components present in a layered network architecture;

Figure 3 is a functional block diagram illustrating the flow of a data packet through program components in accordance with one presently preferred embodiment of the invention;

Figure 4 is a diagram illustrating one presently preferred embodiment of the data packet and the packet extension;

Figure 5 is flow chart illustrating one presently preferred embodiment of the program steps used to offload tasks on per-packet basis; and

Figure 6 is a flow chart illustrating one presently preferred embodiment of the program steps used to query and set the task offload capabilities of a peripheral device.

## 3

### 4

#### 5

## 6 7

## 8

### 9

### 10

## 11

#### 12

The state of the s

12 74

## 13

## 14

## 15

### 16

## 17

## 18

## 19

20

### 21

#### 22

## 23

#### 24

#### DETAILED DESCRIPTION OF THE INVENTION

The invention is described below by using diagrams to illustrate either the structure or processing of embodiments used to implement the system and method of the present invention. Using the diagrams in this manner to present the invention should not be construed as limiting of its scope.

The present invention contemplates both methods and systems for offloading processing tasks from a host computer, such as a personal computer, to computer connected hardware peripherals, such as a network interface card (NIC). Figure 1 and the following discussion are intended to provide a brief, general description of a suitable computing environment in which the invention may be implemented. Although not required, embodiments of the invention will be described in the general context of computer-executable instructions, such as program modules, being executed by a personal computer. Generally, program modules include routines, programs, objects, components, data structures, etc. that perform particular tasks or implement particular abstract data Moreover, those skilled in the art will appreciate that the invention may be practiced with other computer system configurations, including hand-held devices, multiprocessor systems, microprocessor-based or programmable consumer electronics, network PCs, minicomputers, mainframe computers, and the like. Embodiments of the invention may also be practiced in distributed computing environments where tasks are performed by remote processing devices that are linked through a communications network. In a distributed computing environment, program modules may be located in both local and remote memory storage devices.

With reference to Figure 1, an exemplary system for implementing the invention includes a general purpose computing device in the form of a conventional personal

computer 20, including a processing unit 21 (sometimes referred to as the CPU), a system memory 22, and a system bus 23 that couples various system components including the system memory to the processing unit 21. The system bus 23 may be any of several types of bus structures including a memory bus or memory controller, a peripheral bus, and a local bus using any of a variety of bus architectures. The system memory includes read only memory (ROM) 24 and random access memory (RAM) 25. A basic input/output system 26 (BIOS), containing the basic routines that help to transfer information between elements within the personal computer 20, such as during start-up, is stored in ROM 24. The personal computer 20 further includes a variety of peripheral hardware devices, such as hard disk drive 27 for reading from and writing to a hard disk, not shown, a magnetic disk drive 28 for reading from or writing to a removable magnetic disk 29, and an optical disk drive 30 for reading from or writing to removable optical disk 31 such as a CD ROM or other optical media. The hard disk drive 27, magnetic disk drive 28, and optical disk drive 30 are connected to the system bus 23 by a hard disk drive interface 32, a magnetic disk drive-interface 33, and an optical drive interface 34, respectively. The drives and their associated computer-readable media provide nonvolatile storage of computer readable instructions, data structures, program modules and other data for the personal computer 20. Although the exemplary environment described herein employs a hard disk, a removable magnetic disk 29 and a removable optical disk 31, it should be appreciated by those skilled in the art that other types of computer readable media which can store data that is accessible by a computer, such as magnetic cassettes, flash memory cards, digital video disks, Bernoulli cartridges, random access memories (RAMs), read only memories (ROM), and the like, may also be used in the exemplary operating environment.

24

1

2

3

4

5

6

7

8

10

11

12

13

14

15

16

17

18

19

20

21

22

23

1.4. 5 dis

2

3

4

5

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

i, 25 4

1, 5 

ATT OF STATE OF STATE

A number of program modules may be stored on the hard disk, magnetic disk 29, optical disk 31, ROM 24 or RAM 25, including an operating system 35, one or more application programs 36, other program modules 37, and program data 38. A user may enter commands and information into the personal computer 20 through input devices such as a keyboard 40 and pointing device 42. Other input devices (not shown) may include a microphone, joystick, game pad, satellite dish, scanner, or the like. These and other input devices are often connected to the processing unit 21 through a peripheral hardware device such as a serial port interface 46 that is coupled to the system bus, but may be connected by other interfaces, such as a parallel port, game port or a universal serial bus (USB). A monitor 47 or other type of display device is also connected to the system bus 23 via a peripheral hardware interface device, such as a video adapter 48. In addition to the monitor, personal computers typically include other peripheral output devices (not shown), such as speakers and printers.

The personal computer 20 may operate in a networked environment using logical connections to one or more remote computers, such as a remote computer 49. The remote computer 49 may be another personal computer, a server, a router, a network PC, a peer device or other common network node, and typically includes many or all of the elements described above relative to the personal computer 20, although only a memory storage device 50 has been illustrated in Figure 1. The logical connections depicted in Figure 1 include a local area network (LAN) 51 and a wide area network (WAN) 52. Such networking environments are commonplace in offices enterprise-wide computer networks, intranets and the Internet.

When used in a LAN networking environment, the personal computer 20 is connected to the local network 51 through a peripheral hardware device often referred to as

THE REAL PROPERTY.

a network interface card (NIC) or adapter 53. When used in a WAN networking environment, the personal computer 20 typically includes a modem 54 or other means for establishing communications over the wide area network 52, such as the Internet. The modem 54, which may be internal or external, and typically is connected to the system bus 23 via the serial port interface 46. In a networked environment, program modules depicted relative to the personal computer 20, or portions thereof, may be stored in the remote memory storage device 50. It will be appreciated that the network connections shown are exemplary and other means of establishing a communications link between the computers may be used.

Embodiments within the scope of the present invention also include computer readable media having executable instructions. Such computer readable media can be any available media which can be accessed by a general purpose or special purpose computer. By way of example, and not limitation, such computer readable media can comprise physical storage media such as RAM, ROM, EEPROM, CD-ROM or other optical disk storage, magnetic disk storage or other magnetic storage devices, or any other medium which can be used to store the desired executable instructions and which can be accessed by a general purpose or special purpose computer.

When information is transferred or provided over a network or another communications connection (either hardwired, wireless, or a combination of hardwired or wireless) to a computer, the computer properly views the connection as a computer-readable medium. Thus, any such a connection is properly termed a computer-readable medium. Combinations of the above should also be included within the scope of computer readable media. Executable instructions comprise, for example, instructions and data which cause a general purpose computer, special purpose computer, or special purpose

i.

1, 29

The second secon

12.5

ľij.

processing device to perform a certain function or group of functions. Finally, embodiments within the scope of the present invention comprise a computer readable medium having a plurality of data fields stored thereon that represent a data structure.

Embodiments of the present invention are directed to providing the ability to reducing the processing overhead and memory usage of a processing unit 21. This is accomplished by offloading particular computing tasks, which are accomplished for instance by way of an operating system, application programs and/or other program modules that are executing on the processing unit/CPU 21, to an appropriate peripheral hardware device connected to the computer system 20. Many such peripheral devices are increasingly equipped with dedicated processors and memory, and are fully capable of performing many of the same tasks that are typically accomplished solely by the CPU 21. Examples of such devices could include, for instance, network interface cards (53 in Figure 1); disk drive interface cards (e.g., 32, 33, 34 in Figure 1); small computer system interface (SCSI) devices; intelligent serial interface cards; or application specific peripherals, such as devices for the encryption/decryption of data.

While the general inventive concepts discussed herein could be used to offload computing tasks in connection with any of the above peripheral hardware devices, the invention will be described with respect to an example of one presently preferred embodiment, wherein computing tasks are offloaded to a network communications device, such as NIC 53 illustrated in Figure 1. More particularly, illustrative embodiments are discussed as being implemented in the networking environment and architecture of the Windows NT operating system available from Microsoft Corporation. However, it will be appreciated that while specific reference is made to Windows NT concepts and terminology, those skilled in the art will recognize that many, if not most, operating

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

The state of the s

systems and networking architectures share similarities relevant to the environment of the present invention.

In order to more fully understand the context of the present invention, reference is next made to Figure 2, which illustrates a simplified diagram of certain of the components that make up the Windows NT networking model. For purposes of illustration, the OSI layers that correspond to the various Windows NT components are also shown. At the bottom layer, corresponding to the physical layer in the OSI model, resides the actual NICs (sometimes referred to as network cards, or network adapters) 100-104. NICs are the hardware devices that provide the physical interconnection with the physical medium (the network cable), and the transmission of the signals that carry the data generated by all the higher level layers, in accordance with the particular network topology. As noted above, many NICs are equipped with a dedicated processor and memory, and are capable of performing additional sophisticated computing tasks -- including tasks that may otherwise be handled by the host processor CPU. NICs can be physically implemented as a printed circuit board card that is positioned within a slot in the computer, as a PCMCIA type card that is placed within a PCMCIA-compliant slot, as a dedicated chip positioned within the computer chassis on the mother board, or in any other suitable matter.

Each NIC is logically interconnected with the Windows NT networking model, as is schematically represented by bidirectional lines 108-112, via a corresponding network driver 116-120. Network drivers reside in the MAC sublayer of the network model, and link Windows NT to the physical network channel via the corresponding NICs. Each driver, typically implemented as a software component provided by the vendor of the corresponding NIC, is responsible for sending and receiving packets over its corresponding network connection and for managing the NIC on behalf of the operating system. Each ACT IN SECURITY WAS ALSO BEEN ASSESSED. 1 40 ļ, da F 1, 21 A Maria 1

2

3

4

5

6

7

8

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

driver also starts I/O on the corresponding NIC and receives interrupts from them, and calls upward to protocol drivers to notify them of its completion of an outbound data transfer. Also, the device driver will be responsible for invoking, controlling and/or monitoring any of the additional processing capabilities of the corresponding NIC.

In some environments the driver component is written so as to implement a single specific network protocol, such as TCP/IP or XNS. The basic invention of the present invention described and claimed herein would be applicable to such an environment. For purposes of illustration however, the present invention is described in connection with the Windows NT network architecture, in which an interface and environment called the network driver interface specification (NDIS) is provided. The NDIS interface is functionally illustrated in Figure 2 at 126. NDIS shields each of the network drivers 116-120 from the details of various transport protocols (examples of which are shown at 128-134), and vice versa. More particularly, NDIS describes the interface by which one or multiple NIC drivers (116-120) communicate with one or multiple underlying NICs (100-104), one or multiple overlying transport protocol drivers, or transports, (represented at 128-134 in Figure 2), and the operating system.

Essentially, NDIS defines a fully abstracted environment for NIC driver development. Thus, for every external function that a NIC driver needs to perform, from registering and intercepting NIC hardware interrupts to communicating with transport protocol drivers to communicating with an underlying NIC via register manipulation and port I/O, it can rely on NDIS APIs to perform the function. To provide this level of abstraction and resulting portability, NDIS uses an export library referred to as the NDIS Interface Library Wrapper (not shown). All interactions between NIC driver and protocol driver, NIC driver and operating system, and NIC driver and NIC are executed via calls to

The state of the second state of the state o ļadz Harm France was some 1

2

3

4

5

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

wrapper functions. Thus, instead of writing a transport-specific driver for Windows NT, network vendors provide the NDIS interface as the uppermost layer of a single network driver. Doing so allows any protocol driver to direct its network requests to the network card by calling this interface. Thus, a user can communicate over a TCP/IP network and a DLC (or an NWLINK, or DECnet, VINES, NetBEUI and so forth) network using one network card and a single network driver.

At the network and data link layers are transport, protocol and related drivers, shown by way of example in Figure 2 at 128-134. In Windows NT, a transport protocol driver is a software component that implements a transport driver interface (TDI), or possibly another application-specific interface at its upper edge, to provide services to users of the network. In Windows NT, the TDI provides a common interface for networking components that communicate at the Session Layer, such as the Redirector and Server functions illustrated at functional block 138. As is well known, transport protocols act as data organizers for the network, essentially defining how data should be presented to the next receiving layer and packaging the data accordingly. They allocate packets (sometimes referred to in the Windows NT context as NDIS packets), copy data from the sending application into the packet, and send the packets to the lower level device driver by calling NDIS, so that the data can be sent out onto the network via the corresponding NIC.

It will be appreciated that additional functions, or tasks, can also be performed on the data packet as it passes through the various network layers, typically at layers 3 and 4 of the network model. In accordance with the present invention, such additional functions or task may be performed instead by the NIC hardware.

24

WORKMAN, NYDEGGER & SEELEY

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

Hard the transfer and t

For instance, one task conventional performed by transport protocol drivers may be to calculate a checksum value and then append it to the packet. This helps to assure the integrity of the data as it traverses network links. Generally, this operation requires the transport protocol corresponding with the sender of the network packet to append it with a number calculated by adding up the data elements composing the packet. The receiver of the packet then compares the appended checksum number to the data, thereby confirming that data was not changed in transit. This checksum calculation and comparison may instead be offloaded to the NIC hardware.

Another related task that could optimally be performed in the NIC hardware is the calculation of a message digest for the data packet. Like the checksum, a message digest is used to guarantee the integrity of the data in the packet. In addition, a message digest can be used to guarantee the authenticity of the data by assuring that the party who sent the message is who they purport to be. Calculation of a message digest is very CPU intensive, and is a function that is expensive to implement in software.

Another desirable function is the encryption of the data within the packet. Encryption refers to the cryptographic process of transforming the message in the packet so that it becomes impossible for an unauthorized reader of the packet to actually see the contents of the message without prior knowledge of the encryption key. Of course, cryptographic algorithms also tend to be very CPU and memory intensive, and can be prohibitively expensive if performed in software. Examples of such encryption include Secure Socket Layer protocol (SSL) encryption and Internet security protocol encryption or "IPSec." As is well known, SSL and IPSec encryption protocols are both very CPU and memory intensive.

- Page 20 -

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

The party was some of the first

P L J

2 2 2

٠<u>.</u> ٢ 

60 EAST SOUTH TEMPLE SALT LAKE CITY, UTAH 84111

Another task that can be performed on the data packet is TCP or UDP segmentation. As is well known, TCP and UDP protocols segment large data packets into segments that align with the maximum data size allowed by the underlying network. For instance, Ethernet allows a maximum of 1514 byte packets on the network. Thus, if TCP or UDP must send 64 Kbytes for example, it must parse the data into 1514 byte segments. In addition, packet classification is often performed on data packets as they are

received from a network. Packet classification includes the identification of data packet for Quality of Service (QoS). In other words, each data packet contains fields which define modes of service that should be performed on the data packet in order to ensure optimal performance. For example, data packets that include video and audio may require certain functions to be performed thereon in order to ensure high fidelity in the presentation of video and audio. The classification of the data packets in order to identify the mode of service may be offloaded to the NIC hardware in accordance with the present invention.

Packet filtering may also be performed by the NIC hardware instead of by the CPU. In particular, data packets may be evaluated to determine whether they are characteristic of being part of a denial of service attack. Often, server software is created under the assumption that clients that request service are not malicious. However, a client having sufficiently educated, but malicious, user may often issue requests to the server that are designed to prevent the server from servicing other users. Denial of service attack filtering refers to the ability to evaluate data packets to determine if they are characteristic of a denial of service attack.

One example of a denial of service attack is a "SYN flood attack," in which, for one reason or another, the client does not send any final acknowledgment to the server's

The state of

ļ

j" Lji

Maria Maria Maria Maria

SYNchronize-ACKnowledge (SYN-ACK) response in the handshaking sequence. This causes the server to continue signaling until the server eventually times out.

Another type of denial of service attack is called a fragmentation or "teardrop" attack. Internet Protocol requires a packet be fragmented if the packet is too large for handling in the next router. The later fragmented packet identifies an offset to the beginning of the first packet that enables the entire packet to be reassembled at the receive side. In the fragmentation attack, the user puts a counterfeit offset value in the later fragment, often causing the receive side to become dysfunctional. The NIC hardware may be configured to handle filtering for these types of denial of service attacks by evaluating data packets to determine if they are characteristic of a particular kind of denial of service attack.

These and other functions are typically performed by the computer CPU 20 in software components residing at the various network layers, and thus can utilize substantial computer resources, resulting in an overall decrease in the computer system performance. Thus, offloading these, or other similar tasks, so that they can instead be performed at the corresponding NIC hardware can greatly increases the overall speed and efficiency of the computer system.

As previously noted, the basic unit of data transmission in a Windows NT or similar layered networking model is the data packet. In the Windows NT environment, the data packet is referred to as the NDIS packet. Each packet travels from the top of the stack (*i.e.*, layer 5 in the ISO stack) to the lowest software layer (*i.e.*, layer 2 in the ISO stack). Thus, the packet defines a data structure that is common through each level as it proceeds through the layers during transmission and reception of data. By way of example, Figure 3 illustrates the path followed by the packet as it proceeds down through the respective

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

THE THE THE THREE THREE

a de

layers to the NIC, shown at 100 as an Ethernet NIC. As noted above, the transport driver 128 receives data from a sending application and packages it in packet form consistent with the underlying protocol, and then forwards the packet to the lower level device driver 116 via the NDIS interface 126. In addition, the transport protocol may perform other functions on the packet (e.g., checksum calculation, etc.). Alternatively, other functional components may reside in the network layer or data link layers that perform additional functions on the packet, such as the IP Security function 144 (e.g., encryption and/or message digest calculation) illustrated in Figure 3.

In one preferred embodiment of the present invention, the data packet 142 is the means by which computing tasks are offloaded to the peripheral device, such as the NIC hardware 100. For instance, in Figure 3 the application data 140 is passed down from the upper layers of the network model to an appropriate transport protocol driver, such as TCP/IP 128. The driver repackages the data into an appropriate data packet 142. Then, depending on whatever additional functions are to be performed on this particular data packet 142 a functional component is included that appends a predefined data structure, referred to as the packet extension, to the data packet. As will be discussed in further detail below, the contents of the packet extension indicate which task, or tasks, are to be performed on the data packet when it reaches the NIC 100. When the data packet 142 reaches the network driver 116, the contents of this packet extension are queried by the driver 116 so as to ascertain which task(s) is to be performed by the NIC 100. The driver 116 then controls/manipulates the hardware on the NIC so that it will perform whatever functional tasks have been requested via the contents of the packet extension.

For example, in Figure 3, the data packet 142 is passed to a software component 144, which could be implemented separately or implemented as a part of the transport

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

- 1:

protocol driver itself, that appends a packet extension to the packet 142. Data will be included within in packet extension depending on the particular task that is to be offloaded. For instance, if an IP security function is to be implemented, data that indicates that the NIC 100 should encrypt the data packet in accordance with a specified encryption key would be included. Of course, the software component 144 could append predefined data such that any one of a number of functions, such as those discussed above, would be performed at the hardware level instead of by software components that reside in the The device driver 116 will extract the information from the packet network layers. extension, and then invoke the specified task(s) at the NIC 100.

Figure 4 illustrates one presently preferred embodiment of the general structure of the data packet 142. While the packet 142 can be of any format depending on the exact network environment being used, in a Windows NT environment the packet is formatted according to NDIS, and includes information such as a packet descriptor; flags whose meaning is defined by a cooperating device driver(s) and protocol driver(s); areas for storage of Out-of-band data (OOB) associated with the packet; information relating to the length of the packet; and pointers to memory locations relating to the data content of the packet.

Figure 4 further illustrates the additional data structure field that is appended to the NDIS data packet to identify task offloads -- the packet extension 150. As discussed, it is this packet extension 150 which defines a data structure containing information necessary for the identification of the particular task, or tasks, that are being offloaded to the destination NIC. In the preferred embodiment, for each task offload type (e.g., checksum, encryption/decryption, etc) a predefined data field will be included within the packet extension 150. This data field can simply be in the form of a control flag or flags, which

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

ľ.

1-11

merely indicates that a particular function be performed (such as a checksum), or the information can be in the form of a pointer to a data structure that further defines how a task should be carried out.

The NIC 100 may be configured to recognize the task offload control field in the packet extension 150 as applying only to the attached packet. Thus, for example, if a particular packet contains a flag signifying the NIC 100 is to perform a checksum operation, the NIC 100 will perform a checksum operation on the attached packet only. However, if there is no such flag for a given packet, then the NIC 100 does not perform the checksum operation for that packet. Alternatively, the task offload control field in the packet extension 150 may indicate that until the NIC 100 is instructed otherwise, the NIC is to perform the offloaded task on this and all subsequent packets sent out over the network.

The packet extension 150 may also instruct the NIC station at the receive station of what tasks to perform. For example, the packet extension may instruct the receive station NIC that it is to perform certain functions appropriate when receiving a packet from a network such as, for example, decryption, checksum operations, packet classification, packet filtering to guard against denial of service attacks, packet reassembly, and any other receive function that an NIC may be configured to perform. Of course, the sending station may have no idea on the task offload capabilities of the receive station NIC. If the receive station NIC is unable to perform the requested function, then receive station CPU will accomplish the function instead. Thus, the sending station may have some control over what the NIC at the receiving station does to a particular packet.

In addition, the sending station may also use packet extensions to control the NIC 100 when it functions to receive packets. For example, a particular packet to be sent out

24

1

2

3

4

5

6

7

8

9

10

over the network may include a packet extension that includes an instruction to the NIC 100 that when it receives packets, it is to perform decryption of the packet. In this case, the instruction does not apply to the packet to be sent out over the network since the instruction is for performing functions to packets received from the network. Therefore, in this case, the packet extension for a packet to be sent over the network serves as a mechanism to control what receive functions are offloaded to the NIC 100 when receiving packets. The sending station has an opportunity to control the NIC 100 receive functions whenever it sends a packet, on a packet-by-packet basis.

Referring again to the example shown in Figure 4, the packet extension 150 contains a number of data fields 151 to 153 that control the sending functions offloaded to an NIC (e.g., NIC 100) and number of data fields 154 to 158 that control the receiving functions offloaded to an NIC (e.g., NIC 100 or an NIC of a receiving station).

For example, the packet extension 150 includes a data field 151 signifying that the NIC 100 perform a checksum operation. This indicates to the sending NIC 100 that the sending NIC 100 itself is to perform a checksum operation on the appended packet. By way of example and not limitation, a preferred embodiment of the packet extension data structure for signifying that the NIC perform a checksum operation has the following structure:

```
typedef struct NDIS TCP IP CHECKSUM PACKET_INFO
union
      struct
        ULONG
                  NdisPacketChecksumV4:1;
                 NdisPacketChecksumV6:1;
        ULONG
```

```
2
                                                                                                                                                                                                              3
                                                                                                                                                                                                              4
                                                                                                                                                                                                              5
                                                                                                                                                                                                              6
                                                                                                                                                                                                              7
                                                                                                                                                                                                                           8
                                                                                                                                                                                                                     9
THE THE TAX OF THE PARTY OF THE
                                                                                                                                                                                      10
                                                                                                                                                                                11
                                                                                                                                                                                            12
                                                                                                                                                                                            13
      ļa sis
                                                                                                                                                                                            14
      B R B R G""
                                                                                                                                                                                            15
                                                                                                                                                                                            16
                                                                                                                                                                                            17
                                                                                                                                                                                                  18
                                                                                                                                                                                                  19
                                                                                                                                                                                            20
                                                                                                                                                                                            21
                                                                                                                                                                                                  22
                                                                                                                                                                                                  23
```

1

```
Transmit;
     struct
        ULONG
                 NdisPacketTcpChecksumFailed:1;
                 NdisPacketUdpChecksumFailed:1;
        ULONG
                 NdisPacketIpChecksumFailed:1;
        ULONG
                 NdisPacketTcpChecksumSucceeded:1;
        ULONG
                 NdisPacketUdpChecksumSucceeded:1;
        ULONG
                 NdisPacketIpChecksumSucceeded:1;
        ULONG
        ULONG
                 NdisPacketLoopback:1;
         Receive;
     ULONG
               Value;
 };
NDIS TCP IP CHECKSUM PACKET INFO,
*PNDIS TCP IP CHECKSUM PACKET INFO;
```

In this particular packet extension data structure, if the variables NdisPacketChecksumV4 and NdisPacketChecksumV6 both are not set, then the device driver should send the data packet without doing any checksum on it.

The packet extension 150 also may include a security function data field 152, to indicate that the NIC 100 should perform security features such as would be performed in connection with an SSL or IPSec encryption of packet data and/or the calculation of a message digest. For this type of security task, field 152 preferably contains a pointer to a memory location containing a data structure (e.g., data structure 160), which in turn contains information relevant to the performance of the encryption and/or message digest functions. Under some circumstances, the inclusion of a pointer to a memory location having pertinent data has advantages over storing actual data within the packet extension itself.

One such advantage is illustrated in Figure 5, which illustrates an example of a preferred program sequence that corresponds to the situation where multiple successive

1

2

4

5

6

7

8

9

10

11

data packets may need to have the same type of encryption or digest calculation operation performed. After beginning at program step 302, program step 304 determines whether the current packet is the first packet in the sequence of packets. If so, then this first packet is provided with information, or context, that is to be used for that operation on the succeeding packets as well. For instance, the first packet will have a packet extension that sets forth the particular encryption key to use. This value, or context, will be stored in a separate memory location, or handle. Successive packets will not have to include this information, and would, for instance, only have provide a pointer to the memory location where the encryption key (or other context information) is stored, as is denoted at program This approach reduces the overall size of subsequent data packets in the step 308. sequence of packets, and further enhances the efficiency and portability of the task offload method.

The packet extension 150 may also include a segmentation data field 153 (either for TCP or UDP packet segmentation), which instructs the NIC 100 to perform segmentation of bulk data into constituent packets. In this case, the packet extension 150 may be appended to bulk data instead of a packet. The flags in the packet extension may then be seen to apply to all of the resulting packets.

In addition to these send functions performed on packets sent out by the NIC 100, the NIC 100 may be instructed to perform certain offloaded receive functions when receiving packets from the network. These receive function instructions may be given in the packet extension 150 even though the appended packet is sent out over the network, and thus the NIC 100 does not perform the receive function on the appended packet. Examples of such offloadable receive functions are illustrated by the data fields 154

24

23

1 13 1 1000 EAGLE GATE TOWER 60 EAST SOUTH TEMPLE SALT LAKE CITY, UTAH 84111 1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

And the second of the second o

ĽŅ.

PI

e de

1 4:

Some of these offloadable receive functions closely correspond to through 158. offloadable send functions. Others have no close corresponding send function.

Referring to Figure 4, a checksum data field 154 indicates that for all received packets for which checksum calculation has been performed at the send station, the NIC 100 is to perform a checksum calculation to make sure that the data has not been altered en route to the NIC 100. This checksum field 154 may control the checksum receive function for all packets indefinitely or for a given time period unless expressly contravened by a subsequent checksum field.

The packet extension may also include a security flag 155 which indicates that the NIC 100 should perform security features as packets are received such as would be performed in connection with an SSL or IPSec decryption of packet data and/or the calculation of a message digest. For this type of security task, field 155 preferably contains a pointer to a memory location containing a data structure (e.g., data structure 160), which in turn contains information relevant to the performance of the encryption and/or message digest functions.

The packet extension may also include an assembly data field 156 which instructs the NIC 100 to assemble the received packets into bulk data, a classification field 157 which instructs the NIC 100 to classify each packet for quality of service as described above, and a Denial of Service (DOS) attack filter field 158 for filtering incoming packets for characteristics of a DOS attack as described above. The packet extension 150 may also include receive function data fields similar to flags 154-158 intended for the receiving station NIC to perform on the attacked packet.

In a preferred embodiment, the information contained within the packet extension 150 is queried by the particular device driver to which the packet 142 is sent. In the

A STATE OF THE STA

::1

i Li

= ===

Windows NT environment described in the illustrated embodiments, this type of function would preferably be performed by making appropriate NDIS function calls. For instance, a call to a predefined NDIS function that returns a pointer to the packet extension 150 memory location for the packet could be performed. The device driver software could then identify which tasks are to be performed and, depending on the task(s) offloaded, operate/manipulate the driver's corresponding NIC hardware in the appropriate manner.

Utilizing the actual data packet to offload computing tasks from the computer processor to the hardware peripheral is advantageous for a number of reasons. For example, the transport driver can utilize the capabilities of the peripheral on a packet-by-packet basis. This allows tasks to be downloaded dynamically, and the capabilities of a peripheral can be used on an as-needed basis. Thus, if the processing overhead for the computer system is low at a particular point in time, it may be desirable to perform certain tasks on the computer processor in a conventional fashion. Alternatively, if CPU is heavily loaded with other computing tasks, then it can offload tasks to peripheral devices by merely appending the requisite packet extension to the data packets.

Another advantage is the ability offload multiple tasks by way of a single packet, and essentially "batch" a number of operations at once. For instance, when the computer processor performs a checksum operation, or an encryption operation, the entire data field must be loaded into a memory location before the operation can be completed, *i.e.*, the calculation of the checksum or the encryption of the packet data. Moreover, because of the layered networking model only one operation can be performed at a time, thereby requiring the data to be copied into memory multiple times. However, the per-packet approach allows multiple tasks to be offloaded in one packet. Thus, the hardware peripheral can perform two or more operations in a single pass on the data, depending on the capabilities

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

THE STATE WHITE THE REAL PRINCIPLE STATE S

ij

ļeda

ı dı

of the hardware, and thereby greatly increase the throughput and efficiency of the computer system.

It will be appreciated that while the above scheme is particularly useful in connection with the ability to specify tasks to be offloaded to a particular NIC, the perpacket transfer of information could be used in other ways as well. For instance, if a particular NIC is capable of scheduling the delivery of packets at predetermined times, the packet extension data structure could be used to pass information identifying how and/or when the packet should be sent by the NIC hardware.

In a preferred embodiment of the present invention, before a transport protocol driver will append a packet extension to a data packet so as to offload a particular task to a NIC, two additional functions are initially performed. Insofar as there are many different types of hardware peripherals, each having varying processing capabilities, embodiments of the present invention preferably provide a means by which the transport driver can first query the task offload capabilities of the peripherals that are connected to the computer system. Once these capabilities have been ascertained, the transport protocol driver can then set, or enable, those tasks that are of interest. Once enabled, the specified tasks can be subsequently utilized on a per-packet basis in the manner described above.

Figure 6 illustrates one presently preferred set of program steps for implementing this ability to query the task offload capabilities of a peripheral, such as a NIC, and for then setting/enabling those tasks that may be needed. Preferably, the illustrated program steps are implemented as a software component or module that is integrated within a transport protocol driver. Moreover, if implemented in a Windows NT environment, many of the functional operations and inter-driver communications are preferably carried out via the NDIS interface. Of course, if implemented in a different operating system environment, or

implemented in connection with non-network types of peripherals, then the illustrated program steps would have to be modified accordingly.

Beginning at program step 202, executable program instructions are performed by the computer system CPU (e.g. processing unit 21 in Figure 1) that query the task offload capabilities of a peripheral(s) connected to the computer system. In the exemplary environment, the peripheral(s) would be the NIC(s) attached to the system. Preferably, each NIC device driver (116-120 in Figure 2) will have associated with it a predefined task offload buffer location(s) that each contains the task offload capabilities for that device driver and its corresponding NIC. At a minimum, the task offload buffer will identify the particular task(s) supported by the NIC and its device driver, and will also include any information specific to each individual task supported. In the preferred embodiment, the contents of a device driver's task offload buffer are retrieved by way of an NDIS function call.

By way of example, and not limitation, a preferred embodiment of the driver's task offload buffer has the following structure:

```
typedef struct _NDIS_TASK_OFFLOAD

{
ULONG Size;
NDIS_TASK Task;
ULONG OffsetNextTask;
ULONG TaskBufferLength;
UCHAR TaskBuffer[1];

NDIS_TASK_OFFLOAD,
*PNDIS_TASK_OFFLOAD;
```

wherein: Size

}

This variable is set to the size of this structure, and is used to determine the version of the structure in use.

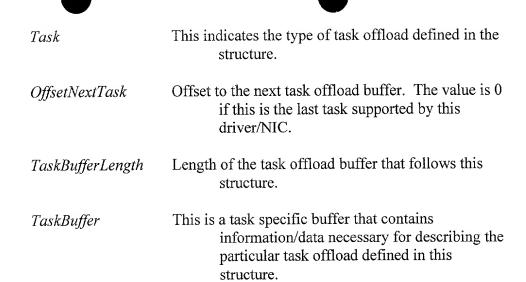
1,13

ĽÜ

1714

-4-

-4-



Once the particular task offload capabilities of each peripheral have been ascertained, processor 21 proceeds to the computer executable instructions corresponding to the function illustrated at program step 204 in Figure 6. Here, it is determined whether any of the particular task offload capabilities supported by each peripheral NIC are of any interest to the transport. If no tasks are supported by the peripheral, or if the tasks that are supported are not useful to the transport, then the processor will proceed directly to program step 208, where this particular line of processing stops, and the transport proceeds in its normal fashion. Alternatively, if one or more of the tasks identified are of interest, then program step 206 is performed. At this step, executable instructions are performed that set those particular task offload capabilities the protocol driver wants enabled. This informs the device driver and its corresponding NIC what type of task offloading it may expect on a per packet basis. In a preferred embodiment, the task offload capability is set by a setting appropriate data values in a data structure that is then passed to the corresponding device driver via an NDIS interface function call. For example, for a particular task, the appropriate bits in the task's task offload buffer structure can be set by the transport, thereby enabling the task for that driver/NIC. The transport can thus enable

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

ANN. Term. 1886 South State State, State,

: 4:

ļ. d.

The little of the control of the con

any number of task offload capabilities of each NIC/NIC driver that it may want to later use.

Once the desired offload capabilities have been enabled for each device driver and its corresponding NIC, the computer system processor 21 proceeds to program step 208, where processing for this particular function ends. At this point, the transport driver can utilize each of the task offload capabilities that have been enabled on a per packet basis in the manner previously described.

In summary, embodiments of the present invention provide distinct advantages over what is currently available in the prior art. Specific processing tasks that would otherwise be performed in a computer system's processor and memory, are instead downloaded to a particular peripheral device, or devices, that are connected to the The computing task is then performed by the peripheral, thereby saving computer system resources for other computing tasks. Especially where the offloaded processing task is CPU and/or memory intensive, this task offload scheme can dramatically increase the overall computing efficiency of the computer system. Advantageously, the disclosed method for offloading computing tasks provides a means by which tasks can be offloaded on a dynamic, as-needed basis. As such, the processor is able to offload tasks in instances where it is busy processing other computing tasks and processor overhead is high. Alternatively, when demand on the computer system's processing resources is lower, the processor may instead perform the task on its own. In addition, multiple tasks can be offloaded in batches to a particular peripheral. Oftentimes, a peripheral will be optimized so as to perform such multiple tasks in a much more efficient manner than the computer processor -- again resulting in improvements to the overall system's efficiency. Finally, a method of the present invention provides a processing scheme by which the

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

Burn dark in the state of the s

z ziz

particular task offload capabilities of a peripheral can be queried, and thereafter selectively enabled. In this way, the computer system can easily discern the capabilities of its various peripheral devices, and take advantage only of those processing capabilities that are, or may be, thereafter needed.

The present invention may be embodied in other specific forms without departing from its spirit or essential characteristics. The described embodiments are to be considered in all respects only as illustrated and not restrictive. The scope of the invention is, therefore, indicated by the appended claims rather than by the foregoing description. All changes which come within the meaning and range of equivalency of the claims are to be embraced within their scope.

What is claimed and desired to be secured by United States Letters Patent is:

- Page 35 -

Docket No. 13768.69.1.1